

1 **Sentence Repetition as a Clinical Marker for Mandarin-Speaking Preschoolers with**
2 **Developmental Language Disorder**

3 Danyang Wang¹, Li Zheng², Yuanyuan Lin^{3,4}, Yiwen Zhang^{3,4}, Li Sheng⁵

4 *¹Department of Communication Sciences and Disorders, University of Delaware, Newark,*

5 *Delaware, USA; ²Nanjing Normal University, Nanjing, Jiangsu, China; ³Developmental and*

6 *Behavioral Pediatrics Department, Shanghai Children's Medical Center, Shanghai Jiao Tong*

7 *University School of Medicine, Shanghai, China; ⁴Ministry of Education-Shanghai Key*

8 *Laboratory of Children's Environmental Health, Xinhua Hospital, Shanghai Jiao Tong*

9 *University School of Medicine, Shanghai, China; ⁵Research Centre for Language, Cognition,*

10 *and Neuroscience & Department of Chinese and Bilingual Studies, The Hong Kong Polytechnic*

11 *University, Hong Kong SAR, China.*

12

13

14 **Address correspondence to** Li Sheng, Research Centre for Language, Cognition, and

15 Neuroscience & Department of Chinese and Bilingual Studies, The Hong Kong Polytechnic

16 University, Hong Kong SAR, China. Email: dr-li.sheng@polyu.edu.hk. Tel: 852-27667446, Fax:

17 852-23340185.

18

19

20 The authors have no relevant conflicts of interest to declare.

21 This project was funded by Humanities and Social Sciences projects of the Chinese Ministry of

22 Education (17YJAZH132) awarded to Li Zheng (PI) and Li Sheng (co-PI), and a Pudong One

23 Hundred Award to Li Sheng.

24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

Abstract

Purpose: Sentence repetition (SR) is believed to be a clinical marker for Developmental Language Disorder (DLD) across many languages. This study explored the potential of a self-designed Mandarin SR task (MSRT) to reflect Mandarin-speaking preschoolers' language ability and to differentiate children with and without DLD in this population. Furthermore, we aimed to compare five scoring systems for evaluating children's MSRT performance. **Method:** In study 1, the MSRT was administered to 59 typically-developing (TD) children aged 3;6 (years; months) to 6;5 in China. The task was examined regarding its ability to correlate with language indices derived from children's narrative samples. In study 2, both a TD and a DLD group were recruited to investigate the task's sensitivity, specificity, and likelihood ratios to distinguish between children with and without DLD. **Results:** Study 1 showed that, using four of the five scoring methods, TD children's performance on the MSRT significantly correlated with all the language measures derived from narratives. Study 2 showed that the MSRT was able to differentiate children with and without DLD. **Conclusion:** The MSRT is a promising tool to reflect language abilities and identify DLD in Mandarin-speaking preschoolers. Based on the current evidence, we recommend that researchers and clinicians select the number of errors in syllable method or the binary method when scoring responses to meet their specific needs.

Keywords: sentence repetition, clinical marker, Mandarin, developmental language disorder

47 **Introduction**

48 Children with Developmental Language Disorder (DLD, also known as Specific
49 Language Impairment), demonstrate significant deficits in talking and/or understanding
50 language, and these deficits are not attributed to any physical or neurological conditions (Bishop
51 et al., 2017). DLD negatively affects children's everyday life and school achievement, including
52 social interaction, literacy, and mathematical thinking (Knox & Conti-Ramsden, 2003; McArthur
53 et al., 2000). Although approximately 7% of children are affected by DLD (Norbury et al., 2016;
54 Tomblin et al., 1997), this disorder is notoriously under-detected in real life (e.g., Jessup et al.,
55 2008; Tomblin et al., 1997). There is a pressing need to develop effective screening methods so
56 that children who are at risk of having DLD can receive further diagnostic assessments which
57 allow them to receive timely support from speech-language pathologists and teachers. A good
58 screening task needs to have high sensitivity to not miss any potential cases of DLD; it also
59 needs to be timesaving so that it can be administered at scale. Sentence repetition (SR, also
60 termed sentence recall, sentence imitation, recalling sentences) tasks are good candidates for this
61 purpose, given its utility in differentiating children with and without DLD across languages (e.g.,
62 Conti-Ramsden et al., 2001; Redmond, 2005; Stokes et al., 2006) and its quick administration
63 and scoring. This paper aims to develop and validate a Mandarin SR task as a measure of
64 language abilities and evaluate its classification accuracy in differentiating Mandarin-speaking
65 preschool children with and without DLD.

66 *Nature of the SR task*

67 Various SR tasks have been developed in a myriad of languages. An SR task involves having
68 speakers listen to auditorily-presented sentences one at a time and repeat each sentence verbatim
69 immediately after presentation. This task is widely recognized as a useful measure of individual

70 differences in speakers' language ability (Polišenská et al., 2015). Baddeley (2000) suggested
71 that speakers rely on their long-term semantic and grammatical knowledge to enable the binding
72 of words into larger sentence-level chunks when performing the SR task. In support of this,
73 Klem et al. (2015) conducted a study with 216 children and found significant correlations
74 between children's SR performance and their knowledge of vocabulary and grammar. The
75 authors suggested that the SR task is a complex language task that reflects the integrity of
76 language processing systems at multiple levels, including lexical and grammatical skills, as well
77 as speech perception and speech production. Similarly, Polišenská et al. (2015) identified the
78 involvement of lexical knowledge and morphosyntax in the successful repetition of sentences
79 and claimed that the SR task reflects speakers' general language ability. Although phonological
80 memory is recruited when performing the SR task (Alloway & Gathercole, 2005), Archibald and
81 Joannis (2009) found that their SR task was more sensitive to deficits in linguistic rather than
82 memory abilities.

83 Not only is the SR task sensitive to individual differences in spoken language ability, it is
84 also a good candidate for identifying children with DLD, as it heavily recruits skills that are
85 known to be weak in this population, such as vocabulary, grammar, and phonological memory
86 (e.g., Bishop et al., 2016; Leonard, 2014; Trauner et al., 1995). Children with DLD start
87 expressing meaning with words 11 months later than their typical peers do (Trauner et al., 1995).
88 Throughout preschool years, they continue to demonstrate deficits in receptive vocabulary
89 (Bishop, 1997; Clarke & Leonard, 1996), expressive vocabulary (Leonard et al., 1999; Thal et
90 al., 1999; Watkins, 1995), and novel word learning (Kan & Windsor, 2010). Early school-age
91 children with DLD also demonstrate word retrieval and semantic processing deficits (Sheng,
92 2014). Grammatical (both morphology and syntax) difficulty is a hallmark deficit in individuals

93 with DLD cross-linguistically (Leonard, 2014). English-speaking children with DLD have
94 difficulties using past tense (e.g., walked) and plural inflections (e.g., ducks) (Joanisse &
95 Seidenbert, 1998), following appropriate word orders (Hansson & Nettelbladt, 1995), producing
96 wh- questions (Van der Lely & Battell, 2003), and using adjuncts (Johnston & Kamhi, 1984).
97 Across different languages, children with DLD show grammatical difficulties that are specific to
98 the ambient language, such as inflectional morphology in Hungarian (Leonard et al., 2009) and
99 aspect markers in Chinese (Fletcher, 2005; Hao et al., 2018). An SR task that capitalizes on these
100 areas of known deficits could therefore act as an effective identification tool of DLD.

101 *Use of Sentence Repetition for Identifying Individuals with DLD*

102 Multiple studies have shown that SR is a clinical marker of DLD, by testing children with
103 and without DLD and investigating the task's classification accuracy values, including sensitivity
104 and specificity. Sensitivity refers to the test's ability to accurately capture individuals with DLD,
105 and specificity reflects the test's ability to accurately identify TD individuals. For diagnostic
106 tasks, Plante and Vance (1994) proposed a guideline which considers sensitivity and specificity
107 values below .80 as unacceptable, values of .80-.89 as acceptable, and values at or over .90 as
108 good.

109 Different SR scoring methods were explored in the literature and were shown to affect
110 the classification accuracy of SR tasks. Commonly used scoring systems include (1)
111 correct/incorrect scoring method (binary method; Newcomer & Hammill, 2019; Rispen, 2004),
112 which gives a score of one to completely accurate repetitions and a score of zero to responses
113 with any deviations; (2) scoring that considers errors (error scoring method), in which the score
114 (0-2 in Redmond, 2005; 0-3 in the *Clinical Evaluation of Language Functions*, CELF, Wiig et
115 al., 2013) is based on the number of errors in the response; (3) scoring that considers the specific

116 grammatical structures in the sentences (known as the core element scoring or grammatical
117 scoring, Komeili & Marshall, 2013), which gives a score of one to the response that contains the
118 target grammatical structure and a score of zero to a response without the target structure; (4)
119 scoring that calculates the percent of correct syllables (correct syllable/total syllables; Stokes et
120 al., 2006).

121 Research evidence has supported the use of SR tasks as clinical markers for DLD in
122 English. Conti-Ramsden et al. (2001) compared four potential psycholinguistic markers for DLD
123 in English-speaking children: an SR task (CELF-R-Recalling Sentences subtest, Semel et al.,
124 1994), a nonword repetition task (in which children repeat nonsense words of varied lengths), a
125 third person singular task, and a past tense task. Compared to the other three markers, the SR
126 task demonstrated the highest diagnostic accuracy (using the error scoring), with a sensitivity of
127 90% and a specificity of 85% (using a cutoff score of -1SD). Using the error method, another
128 English SR task developed by Redmond (2005) also demonstrated good utility (sensitivity of
129 94% and specificity of 88%) in differentiating children with and without DLD.

130 SR tasks were shown to be effective clinical markers for DLD in other languages as well.
131 Leclercq et al. (2014) compared group performance in children with and without DLD on a
132 French SR task under seven different scoring methods¹. Children with DLD performed
133 significantly lower than TD children on the French SR task, regardless of the scoring method.
134 When using the correct/incorrect scoring method, the discriminant function analysis revealed the
135 highest levels of sensitivity (97%) and specificity (88%). Furthermore, Armon-Lotem and Meir

¹ The seven scoring methods included the binary scoring; grammatical scoring; scoring that considers three core semantic ideas in each sentence; as well as scorings that calculate the number of correct words; number of correct morphemes; number of function words; and number of content words.

136 (2016) established the effectiveness of their SR tasks in distinguishing monolingual children with
137 and without DLD in both Russian (sensitivity of 86%, specificity of 90%) and Hebrew
138 (sensitivity of 100%, specificity of 87%), although the scoring method was not specified.

139 Only a few studies have investigated the utility of SR as a clinical marker for DLD in
140 Asian languages, including Korean (Hwang, 2012), Vietnamese (Vân Hoàng et al., 2014; Pham
141 & Ebert, 2020), and Cantonese (Stokes et al., 2006). Pham and Ebert (2020) tested 104
142 Vietnamese-speaking five- and six-year-old children, including ten children with DLD. The
143 authors explored three scoring methods: binary, error, and grammatical scoring. When using the
144 error scoring method, their SR task achieved a sensitivity of 90% (CI²: 0.71 to 1.09³) and a
145 specificity of 71% (CI: 0.43 to 0.99). The binary method yielded slightly higher sensitivity
146 (100%) but lower specificity (57%, CI: 0.26 to 0.88) than the error method. When using the
147 grammatical scoring method, the sensitivity was 80% (CI: 0.55 to 1.05) and the specificity was
148 71% (CI: 0.43 to 0.99).

149 Stokes et al. (2006) examined a Cantonese SR task with a DLD group (N=14), an age-
150 matched TD group (N=15) and a younger language-matched TD group (N=15). Four different
151 scoring methods were explored: binary, error, grammatical and percent of correct syllable
152 scoring. The age-matched TD group performed significantly higher than the language-matched
153 TD group and the DLD group, whereas the latter two groups did not differ from each other,
154 regardless of the scoring method. Classification accuracy was evaluated using the error method
155 and the percent correct syllable method. The error method resulted in higher classification values
156 (sensitivity of 77%, CI: 0.55 to 0.99; specificity of 100%) in differentiating between the DLD

² CI=95% confidence interval

³ Although the confidence interval mathematically exceeds 1, the maximum value of sensitivity/specificity is 1.

157 group and the age-matched TD group, compared to the percent correct syllable method
158 (sensitivity of 43%, CI: 0.17 to 0.69; specificity of 100%).

159 To summarize, the SR task has received support as a clinical marker for DLD across
160 languages. How SR responses are scored has a direct bearing on the classification accuracy of
161 the SR task (Leclercq et al., 2014; Pham & Ebert, 2020; Stokes et al., 2006). The error method is
162 the most commonly used scoring system and consistently shows fair to good classification
163 performance in different languages (Pham & Ebert, 2020; Redmond, 2005; Stokes et al., 2006;
164 Wiig et al., 2013). The binary method is the simplest and yielded the highest sensitivity in
165 Vietnamese (Pham & Ebert, 2020) and the best overall classification results in French (Leclercq
166 et al., 2014). For a newly established SR task, these results highlighted the importance of
167 empirically testing the most effective scoring method.

168 *A Sentence Repetition Task in Mandarin*

169 To the best of our knowledge, the utility of the SR task in Mandarin is yet to be
170 established. The investigation of the SR task in Mandarin, a typologically distinct language from
171 Indo-European languages, can further support this task's utility to differentiate individuals with
172 and without DLD across languages and contribute to the understanding of the underlying deficits
173 of DLD (Pham & Ebert, 2020). In addition, Chinese has 873 million native speakers and 178
174 million second language speakers all over the world, and Mandarin speakers constitute the
175 majority of this population (Gordon, 2005). Following the logic that DLD is affected by genetic
176 components and should be equally prevalent across languages and countries (Armon-Lotem et
177 al., 2015; Rice, 2013), there are approximately 5 million 4-9 years old children in China that are
178 estimated to have DLD (Sheng et al., 2020). The development of a Mandarin SR task as a

179 screening tool is thus clinically significant to facilitate the early identification of DLD in
180 Mandarin-speaking children.

181 There is a substantial literature on the manifestations of DLD in Mandarin in the
182 preschool to early elementary school age period (Sheng et al., in preparation) that could guide
183 the design of stimuli for a novel SR task. Specifically, this literature highlights several structures
184 that are well-established in typically-developing children but present considerable challenges to
185 children with DLD, including passives, classifiers, and aspect markers. Classifiers in Mandarin
186 modify nouns that share the same properties in terms of shape or other dimensions (Lin & Bever,
187 2010), and they are mandatory when adding numerals to nouns (e.g., san1 *zhil* gou3; three
188 *classifier* dog). Passive sentences follow non-canonical word order and individuals with DLD
189 across languages exhibit difficulty in this structure (Leonard, 2014). Zeng et al. (2018) found that
190 children with DLD performed significantly lower than their TD age-matched peers on the
191 comprehension and production of Mandarin passives. A narrative study also showed that
192 Mandarin-speaking children with DLD produced significantly fewer passive sentences compared
193 to their TD peers (Hao et al., 2018). In addition, Hao et al. (2018) observed significantly fewer
194 classifiers and aspect markers used by children with DLD than TD children. Weaknesses in
195 aspect markers were also shown in He and Sun (2013), which found that Mandarin-speaking
196 children with DLD performed significantly worse than age-matched TD children on an aspect
197 marker production task.

198 *The Current Study*

199 The current paper reports two studies that respectively examined the concurrent criterion
200 validity of the self-designed Mandarin Sentence Repetition Task (MSRT) against criterion
201 language measures of narrative sampling (study 1), and the discriminant validity (i.e.,

202 classification accuracy) of the MSRT against clinical diagnosis of DLD based on pediatrician
203 judgment and standardized test scores (study 2). In study 1, to establish that the MSRT can
204 reflect children's language ability, we examined if MSRT scores were correlated with measures
205 derived from children's narrative samples. Through collecting and examining functional
206 language use at the discourse level, researchers and clinicians could gain valuable insights into
207 an individual's language abilities in everyday communication (Spencer et al., 2020). The
208 analysis of language samples yields a deep and comprehensive evaluation of children's
209 knowledge in different linguistic domains, including syntax, vocabulary, and use of specific
210 linguistic structures (e.g., Andreu et al., 2011; Boudreau, 2008). Establishing that this newly
211 designed MSRT aligns with measures derived from narrative samples in TD children is a crucial
212 step before moving forward to examine the classification accuracy values in study 2.

213 In study 1, we derived four different measures from narrative samples to evaluate their
214 relationship with children's SR performance, including mean length of utterance (MLU),
215 vocabulary diversity (VOCD), number of predicates, and a composite structural measure. Both
216 MLU and VOCD are commonly used in language sample analysis, as general measures of
217 grammar (Boudreau, 2008; Justice et al., 2010) and vocabulary (Altman et al., 2016; Rezzonico
218 et al., 2015) respectively. We did not use the type-token ratio (TTR) measure because it is
219 subject to the influence of sample length: longer samples may give lower TTR values (Richards,
220 1987; Tweedie & Baayen, 1998). As opposed to a single value of TTR, VOCD is more
221 informative because it represents how TTR varies over a range of token size for each speaker
222 (Richards & Malvern, 2000).

223 Number of predicates is another measure of syntactic elements (Eisenberg, 2020), which
224 include both verbs and predicate adjectives in Mandarin (Thomson & Tao, 2010). Devoscovi and

225 Cristina Caselli (2007) found significant correlations between children’s performance on an
226 Italian SR task and the number of predicates (verbs only in Italian) children used in spontaneous
227 language samples. Moreover, to examine whether production of specific linguistic elements in an
228 imitation context is related to production of the same elements in a spontaneous context, we
229 included a composite structural measure that evaluates children’s use of the linguistic structures
230 that are featured in the SR stimuli, which will be named in the methods section.

231 As scoring methods may directly impact the classification accuracy of SR tasks, we
232 explored five commonly used scoring methods in study 1. The scoring system(s) that did not lead
233 to ceiling/floor effects and demonstrated significant correlations with all narrative measures were
234 retained in study 2. We recruited Mandarin-speaking children with and without DLD in study 2
235 and examined the classification accuracy of the self-designed MSRT to differentiate these two
236 groups. Showing correlations with gold standard measures is insufficient in demonstrating the
237 clinical utility of a task, as reflecting general language abilities in TD children does not equal to
238 accurately identifying DLD on a child-by-child basis. We calculated sensitivity, specificity, and
239 likelihood ratios to evaluate the task’s classification accuracy. We used Receiver Operating
240 Characteristic (ROC) analysis to generate the optimal cutoff scores for the practical use of
241 MSRT as a clinical screening tool. Overall, we aimed to answer three research questions. Study 1
242 addressed the first two questions and study 2 addressed the last question.

243 1) Does the MSRT reflect TD Mandarin-speaking children’s language ability, by showing
244 significant correlations with four narrative measures: MLU, VOCD, number of
245 predicates, and a composite structural measure?

246 2) Out of the five scoring methods, which method(s) are not subject to ceiling/floor effects
247 and can reflect children's language competence by demonstrating significant
248 relationships with the four narrative measures?

249 3) Is the MSRT able to differentiate Mandarin-speaking children with and without DLD?

250 **Study 1**

251 **Method**

252 **Participants**

253 Fifty-nine Mandarin-speaking preschoolers (30 males, 29 females) participated in study
254 1. All participants were Asians of Chinese ethnicity. Parents of the participating children signed
255 an informed consent approved by the University of Delaware's Institutional Review Board.
256 Children's age ranged from 45 to 77 months, and the mean age was 61.7 months. The
257 participants were recruited from the same preschool in Nanjing, China. All children were
258 typically-developing with no reported language, sensory, speech production, motor, or cognitive
259 disorders according to parent reports. Children's nonverbal intelligence was measured using the
260 *Primary Test of Nonverbal Intelligence* (PTONI) (Ehrler & McGhee, 2008), and the average
261 standardized score was 125.1 (SD=17.2, range=85-149). Participants' caregivers filled out a
262 questionnaire to report children's family background, general health condition and
263 developmental history, and Mandarin exposure.

264 Children's family SES was collected through surveying their maternal education, using a
265 five-point likert scale. Twenty-five percent of the parents had a master's degree or higher, 56%
266 had a bachelor's degree, 17% completed some college, and 2% completed middle school or
267 lower. The mean maternal education score was 4 (bachelor's degree) (SD=0.77). Participants'

268 Mandarin exposure was collected using a five-point likert scale question asking about the percent
269 of the waking hours that the child spent hearing and speaking Mandarin. Children's average
270 Mandarin exposure score was 4.63 (SD=0.75), with a score of 4 indicating 60-79% and a 5
271 indicating 80-100%. Other than Mandarin, children were either exposed to dialects that are
272 mutually intelligible with Mandarin (e.g., Henan and Nanjing dialects) and/or were exposed to
273 other languages such as English (n = 40) and German (n = 1). It is worth noting that previous
274 studies of Mandarin-speaking children's language development rarely if ever reported children's
275 Mandarin exposure (Sheng et al., in preparation). Researchers likely assumed the monolingual
276 status of their sample because of the prestige of Mandarin in the Chinese society: Mandarin is the
277 only official language, the language of media, and the language of instruction at schools (Dong,
278 2010). Though the average amount of Mandarin exposure of the current sample seems low, given
279 the similarity in recruitment approaches, we believe the sample is comparable to samples in
280 previous studies of Mandarin language development, and is representative of the language
281 exposure patterns of preschool age children in mainland China.

282 **Test Materials**

283 *Sentence Repetition Task*

284 In keeping with previous studies (Redmond, 2005; Stokes et al., 2006) and the review on
285 the manifestations of DLD in Mandarin, we included eight sentences with passives and eight
286 sentences with aspect markers as core elements. Each sentence type comprises eight test
287 sentences with length ranging from 13 to 15 characters. Ninety percent of the nouns, verbs, and
288 adjectives included in the sentences are early acquired lexical items, which appear in children's
289 production as early as 17 months of age according to the Chinese Communication Development
290 Inventory (CDI; Tardif & Fletcher, 2008). More challenging but age-appropriate words (ten

291 percent, e.g., 美味, mei3-wei4, delicious; 批评, pi1-ping2, criticize; 偷走, tou1-zou3, steal) were
 292 included as well to increase the task's ability of revealing individual difference and avoid ceiling
 293 effect in TD children. The passive sentences constitute the elements of patient, "bei4" (passive
 294 marker), agent, and verb (see 1 for an example). The aspect marker sentences contain subject,
 295 verb, aspect marker, and object (see 2 for an example). Eleven sentences include classifiers
 296 (bolded in Appendix A), and four sentences contain embedded relative clauses (underlined in
 297 Appendix A). Relative clauses are featured in the stimuli to avoid ceiling effect in TD children,
 298 as they are structurally complex and acquired late developmentally (He et al., 2017; Sung et al.,
 299 2016). All sentences were pre-recorded by a female native speaker of Mandarin Chinese.

300 1) 那只白色的小狗被妈妈抱走了。

301 Na4 zhi1 bai2 se4 de xiao3 gou3 bei4 ma1 ma1 bao4 zou3 le.

302 That CL white LP dog PASS mother carry away(RP) SFP

303 That white dog was carried away by the mother.

304 2) 小老鼠吃了一块美味的巧克力。

305 Xiao3 lao3 shu3 chi1 le yi1 kuai4 mei3 wei4 de qiao3 ke1 li4.

306 Little mouse ate PerM one CL delicious LP chocolate.

307 The little mouse ate a piece of delicious chocolate.

308 *CL = classifier; SFP = sentence final particle; PerM = perfective aspect marker; PASS =
 309 passive marker; LP = linking particle; RP = resultative particle

310 *Narrative Task*

311 Children completed the Mandarin version of the Multilingual Assessment Instrument for
312 Narratives (MAIN, Gagarina et al., 2012; Luo et al., 2020). The MAIN was designed based on a
313 series of pilot studies with more than 500 monolingual and bilingual children in 17 languages
314 and 14 language pairs, and the stimulus pictures and scripts were carefully constructed to elicit
315 narrative samples from children with diverse cultural, linguistic, and socio-economic
316 backgrounds (Gagarina et al., 2012). The MAIN encompasses four stories, with two stories
317 assigned to the story-telling format (*dog* and *cat*) and two stories assigned to the story-retelling
318 format (*baby bird* and *baby goat*). All stories are parallel regarding cognitive and linguistic
319 complexity, cultural appropriateness, and test robustness (Gagarina et al., 2012). Each story has a
320 setting and episode structures that can capture the universal organizational pattern of stories. The
321 MAIN has been used in 15 different languages as a language assessment tool (e.g., Dutch: Blom
322 et al., 2020; Greek: Tsimpli et al., 2020; Italian: Levorato & Roch, 2020; Cantonese: Chan et al.,
323 2020; Mandarin: Sheng et al., 2020).

324 **Procedures**

325 *Sentence Repetition Task*

326 Each child was assessed individually in a private room at their school. The task was
327 administered using a PowerPoint presentation on a computer. The experimenter sat at a table
328 next to the child and started with playing the instructions on the computer. The instruction was:
329 “Let’s play an imitation game on the computer. The computer will say some sentences. Your job
330 is to listen carefully to what the computer says and repeat the sentences. Please remember, you
331 have to say exactly the same thing as the computer. Let’s practice!”. Following the instructions,
332 two practice items were administered one at a time to ensure that children understood the task.
333 For the practice items, the experimenter guided the child to repeat the sentences and corrected

334 them if they did not repeat verbatim. The test phase began once the child had successfully
335 repeated the practice sentences word for word. A total of 16 test sentences were then presented
336 one at a time at 65dB SPL using the computer's built-in speaker. The experimenter pressed a
337 button on the computer to play the next sentence once the child had responded. The experimenter
338 was instructed not to interfere with children's performance in any means besides providing
339 general encouragement for children to continue during the actual test. The SR sessions were
340 recorded using a Philips VTR5100 voice recorder with the noise reduction function.
341 Transcriptions were completed later based on the recordings.

342 *Narrative Task*

343 Two tasks were administered in this study: first a story-tell task and then a story-retell
344 task. The child was presented with three envelopes containing the same story inside (*baby goat*
345 or *baby bird*) and was told that each envelope contained a different story. The experimenter
346 pretended that they did not know which story the child would choose to create a more interactive
347 environment and motivate the child to tell the story in detail. After the child had made their
348 decision, the experimenter held the pictures facing the child to give them an overview of the
349 story. When the child was ready to tell the story, the experimenter presented them with two
350 pictures at a time. At the end of each two pictures, the experimenter encouraged the child to say
351 more using prompts like "What else?", "Can you tell me more?", or "Is that all?". In the retell
352 task, the child was presented with another three envelopes containing the same story (*dog* or *cat*).
353 After the child made their decision, the examiner told a model story using the script provided by
354 Luo et al. (2020). Subsequent to the demonstration, the child was asked to retell the story with
355 the aid of the pictures. They were presented with two pictures at a time to ensure that they
356 followed the story in sequence. The same prompts were given at the end of each two pictures.

357 The narrative sessions were recorded using the same voice recorder. Transcriptions were
358 completed later.

359 **Transcription and Scoring**

360 *Sentence Repetition Task*

361 Responses to each sentence were transcribed verbatim using a Microsoft excel worksheet.
362 We explored five potential scoring methods based on the literature. The first two followed the
363 error method (Conti-Ramsden et al., 2001; Stokes et al., 2006; Wiig et al., 2013) and are
364 respectively referred to as error method-word and error method-syllable. In Mandarin, a syllable
365 represents one character, and a word may constitute one (e.g., 我, wo3, “me”) or more
366 syllables/characters (e.g., 公园, gong1-yuan2, “park”). Two error methods are explored because
367 words and syllables are both potential basic units of grammar in Mandarin (Duanmu, 2016) and
368 the counting unit (word vs. syllable) may alter the number of errors. For example, a substitution
369 of “蛋糕, dan4-gao1, cake” with “饼干, bing3-gan1, cookie” would be counted as one error in
370 error method-word and two errors in error method-syllable. In both error methods, a score of
371 three is given to completely accurate repetitions, a score of two is given if one error is found in a
372 response, a score of one is given to a response that contained two to three errors, and a score of
373 zero is assigned to a response that contained four errors or more. Any deviation in words or
374 syllables from the target sentence (e.g., substitution, deletion, addition) was counted as one error.
375 In both error methods, the score for each child was calculated by dividing the child’s score by the
376 total possible score of 48 (3 x 16 sentences). Mazes (e.g., filled pauses, revisions) were not
377 included as errors in the two error methods.

378 The third method is the binary scoring method (Newcomer & Hammill, 2019; Rispens,
379 2004), in which only a completely correct repetition is given a score of 1. Responses with any
380 deviations from the target were scored as 0. The percentage of completely correct repetitions (out
381 of 16) were calculated for each child. The fourth method is the core element method, also known
382 as grammatical scoring (Pham & Ebert, 2020; Stokes et al., 2006). Using this method, a response
383 receives a score of 1 if it contains all the core elements of the two stimulus types. For the passive
384 sentences, the core elements are bei4 + agent + verb, whereas for the aspect marker sentences,
385 the core elements are verb + aspect marker + noun. The last method is the percent of correct
386 syllable method (Stokes et al., 2006), which calculates the percentage of correct syllables.
387 Additions and transpositions were not penalized. As speech production errors may act as a
388 confound in a language task and negatively affect children's scores, we did not take any points
389 off for clear speech production errors (e.g., pronouncing /t^h/ as /t/; 糖, tang2 -> dang2). The
390 illustration of the five scoring methods for an example response is presented in the
391 supplementary material.

392 *Narrative Task*

393 Children's narrative samples were transcribed into Chinese characters using the Codes
394 for Human Analysis of Transcripts (CHAT; MacWhinney, B., 2000). The transcriptions were
395 analyzed using the Computerized Language Analysis (CLAN; MacWhinney, B., 2000). The
396 following measures were calculated from the transcriptions of both the tell and retell samples:

397 1) Mean Length of Utterance (MLU) and vocabulary diversity (VOCD)

398 All transcriptions were segmented into independent clauses which constitute the basis of
399 utterance measures (Sheng et al., 2020). The utterances were first segmented into words

400 using the “*Chinese online word segmentation system*” (<http://ckipsvr.iis.sinica.edu.tw>)
401 and then checked manually. MLU-word and VOCD were generated using the CLAN
402 system.

403 2) Number of predicates: the total number of verbs and predicate adjectives in children’s
404 telling and retelling samples were manually coded and calculated. Verbs include both
405 action verbs (e.g., 吃, *chi1*, eat; 看, *kan4*, look; 跑, *pao3*, run) and modal auxiliary verbs
406 (e.g., 可以, *ke2-yi3*, can; 要, *yao4*, want/will). Predicate adjectives are the adjectives
407 used in sentences without verbs.

408 3) Structural composite: Children’s correct use of classifiers, aspect markers, passives and
409 relative clauses was manually coded and calculated. The general classifier *ge4* was not
410 included in the count of classifiers. Every unique combination of classifier + noun was
411 counted. For example, 一只羊 (one *CL-zhi1* sheep) and 一只鸟 (one *CL-zhi1* bird) were
412 counted as two different classifier uses. Two types of aspect markers were coded and
413 counted in this measure: the progressive markers (*zai4* and *zhe*) and the perfective
414 markers (*le* and *guo4*). There are two different uses of *le* in Chinese: one is a genuine
415 perfective marker indicating the perfective aspect, and the other is a sentence-final
416 particle that marks the reported event or situation as “relevant” to the context (Li &
417 Thompson, 1989; Wang & Sun, 2015). In this analysis, only the correct use of *le* as a
418 perfective aspect marker was counted. Total number of grammatical passive and relative
419 clause sentences was calculated as well. The structural composite score was derived by
420 adding up the number of correct uses of the four structures in children’s narrative samples
421 (tell and retell).

422 **Reliability**

446 VOCD, the structural composite, and number of predicates were normally distributed. SR scores
447 using other scoring methods were not normally distributed. Despite the fact that some variables
448 were normally distributed, we employed the non-parametric Spearman's rank order correlation
449 for all pairs to allow reasonable comparison among the correlational outcomes for the five
450 scoring methods. Though outliers were present, the Spearman's rank order correlation does not
451 require outliers to be removed.

452 Participants' SR scores are presented in Figure 1. Both the core element scoring and the
453 percent of correct syllable scoring led to a ceiling effect. Using the core element scoring, over
454 half of the children (N=33) scored 100% on the MSRT; and using the percent of correct syllable
455 scoring, over 70% of the children (N=45) scored between 90% and 100%. Table 1 shows the
456 results of the Spearman rank correlations. The scatterplots of the relationship between children's
457 SR scores and performance on the narrative measures were presented in the supplementary
458 material. The p level was corrected for multiple correlations using the Bonferroni's correction
459 (corrected significant level: $p < .003$). Four of the five scoring methods demonstrated significant
460 correlations with all narrative measures. The core element scoring did not correlate with MLU,
461 and the number of predicates in narratives.

462 Insert Figure 1 about here

463 Insert Table 1 about here

464 **Interim Discussion**

465 In study 1, we validated the self-designed Mandarin SR task (MSRT) against four
466 language measures derived from children's narrative samples (MLU, VOCD, number of
467 predicates, structural composite). The results demonstrated that using four out of the five scoring

468 methods, children's SR performance significantly correlated with all narrative measures,
469 indicating that the MSRT can reflect Mandarin-speaking children's language abilities. The
470 logical next step is to further evaluate the classification accuracy of this task in differentiating
471 between Mandarin-speaking children with and without DLD.

472 Study 1 results could further guide the selection of scoring methods in the classification
473 accuracy study. Specifically, the core element scoring could be excluded as it correlated with
474 only two of the four narrative measures and showed a ceiling effect. Although the percent of
475 correct syllable scoring did show significant correlations with the narrative measures, it also led
476 to a ceiling effect, which is not desirable in test development. The two error methods and the
477 binary method appeared to be superior to the other methods because they showed significant
478 correlations with narrative measures, and they elicited a wide range of performance. The two
479 error methods showed comparable results but differed in ease of use. The judgment of syllables
480 in Mandarin is more straightforward than that of words, as each syllable is equivalent to one
481 character and each word may contain one or more syllables. The definition of word in Mandarin
482 can be controversial (Li & Thompson, 1989). For example, Jespersen (1922) concluded that
483 Mandarin words are essentially monosyllabic, while Kennedy (1951) and Lin (1952) argued that
484 most Chinese words occur in disyllabic forms. Other researchers even suggest that wordhood is
485 nonexistent in Chinese languages (Chao & Yang, 1947; Chéng, 2003). Accurate word
486 segmentation requires a combination of automated segmentation software and manual correction,
487 which adds considerable time to the scoring process and requires substantial linguistic
488 knowledge from researchers and clinicians. We therefore further explored the error method in
489 syllable and the binary method in study 2.

490

Study 2

491

Method**492 Participants**

493 Sixty-nine Mandarin-speaking children between the ages of 4;0 and 5;11 were recruited
494 and tested. All but one participants were Asians of Chinese ethnicity. One child was of mixed
495 Chinese and Japanese ethnicity. Parents of the participants signed an informed consent approved
496 by the Research Ethics Board of the Shanghai Children’s Medical Center. Participant recruitment
497 was conducted in three phases. In phase one, a one-gate design was applied, with both groups
498 recruited simultaneously from the same site, the Developmental and Behavioral Pediatrics
499 Department at the Shanghai Children’s Medical Center. Thirty-eight children who visited the
500 clinic and were screened as having no physical or neurological impairments participated from
501 March 2019 to December 2019. Further classification of this group yielded five TD children in
502 phase one. Given the inefficiency of recruiting TD children from an outpatient clinic setting, in
503 phase 2, we recruited TD children from local communities using word of mouth and advertising
504 on a popular social media platform in China, WeChat. Ten TD children were recruited and tested
505 in two weeks. Phase 2 was paused in January 2020 because of the outbreak of the Coronavirus
506 Disease in China. Recruitment for the TD group resumed in June 2020, and 21 children were
507 recruited and tested within a month.

508 Participating children first received a hearing screening from a pediatrician using a
509 portable audiometer. Children who have normal hearing then received a clinical screening
510 conducted by a developmental and behavioral pediatrician with over 30 years of experience to
511 rule out other physical and neurological impairments. The decisions were made based on the
512 standards described in the *Diagnostic and Statistical Manual of Mental Disorders-V* (DSM-V;
513 APA, 2013) and children’s medical history. Children who had normal hearing, normal or

514 corrected-to-normal vision, and no diagnosis of neurological disorder, speech production
515 deficits, autism, genetic disorder, or cerebral palsy were invited to complete the *Wechsler*
516 *Preschool and Primary Scale of Intelligence* – Revised (WPPSI-R; Wechsler, 1989). The
517 WPPSI-R is a standardized intelligence scale designed for children aged 3;0 to 7;3. A Chinese
518 version is available (Chen & Chen, 2000), which was normed on 900 Taiwanese children. To be
519 included in study 2, children needed to score over 80 on the WPPSI-R performance scale (.89 of
520 Cronbach’s alpha and test-retest reliability; Chen & Chen, 2000), to demonstrate normal
521 nonverbal intelligence. Four children were excluded from the study because they did not meet
522 the inclusion criterion for performance IQ.

523 Qualified children were then invited to complete the *Diagnostic Receptive and*
524 *Expressive Assessment of Mandarin* (DREAM; Ning et al., 2014), which was used for language
525 status classification. The DREAM is a standardized, norm-referenced oral language assessment
526 for Mandarin-speaking children ages 2;6 to 7;11. The test provides one total score and four
527 component scores: expressive language, receptive language, semantics, and syntax. DREAM
528 achieved high test-retest reliability ($r=.85$) and good external validity by demonstrating
529 significant correlations with spontaneous language indices (e.g., sentence complexity and
530 vocabulary diversity) and narrative indices (e.g., use of mental verbs and connectives) (Liu et al.,
531 2017). When validated against a combination of pediatricians’ judgment and spontaneous
532 language samples, a cutoff score of 80 on any one of the DREAM components yielded a
533 sensitivity of 95% and a specificity of 82% in differentiating children with and without DLD
534 (Liu et al., 2017). We followed these empirically-derived guidelines regarding the use of the
535 DREAM test scores and included children who had at least one component standard score at or
536 below 80 (the 10th percentile or 1.3 SD below the mean) in the DLD group.

537 Of the 69 children recruited in study 2, 27 children met the criteria as having DLD. One
538 child with comorbid ADHD was included in the DLD group, as the most recent DLD definition
539 does not exclude children with ADHD (Bishop et al., 2017). Eight children were excluded as
540 their parents did not complete the parent questionnaire to report family background and language
541 environment. One child was excluded as one of his/her parents is a native Japanese speaker. We
542 then attempted to select an age-matched TD peer for each of the remaining 18 children with
543 DLD. An eligible TD match needed to 1) have all DREAM component scores higher than 80; 2)
544 be within six months of age of the child with DLD; 3) have a Mandarin exposure score within ± 1
545 from the child with DLD; and 4) have a maternal education score within ± 1 from the child with
546 DLD. We were able to find TD matches for 16 children with DLD. As shown in Table 2, the TD
547 and DLD groups did not differ significantly on age, maternal education, and Mandarin exposure.
548 The DLD group showed significantly lower nonverbal IQ scores and DREAM total scores than
549 the TD group.

550 [Insert Table 2 about here](#)

551 **Procedures**

552 Each child participated in the study in a quiet assessment room in the Developmental and
553 Behavioral Pediatrics Department of the Shanghai Children's Medical Center. A trained native
554 Mandarin-speaking research assistant administered the MSRT to the child as part of a larger
555 battery of language and cognitive tests. The MSRT was given in the middle of a larger language
556 assessment battery and the administration followed the same procedures as described in study 1.
557 The error method in syllable and the binary method were used to score children's responses.

558 **Reliability**

559 The first author of this paper transcribed and scored all participants' responses. A second
560 native Mandarin-speaking trained research assistant transcribed and scored 20% of the data
561 independently to examine reliability. During transcription and scoring, both the first author and
562 the research assistant were blinded to the grouping status of the children to avoid potential
563 biases. Transcription reliability was calculated by dividing the number of consistent characters
564 by the number of total characters in the sentences, yielding an inter-rater reliability of 96%.
565 Scoring reliability was calculated by dividing the number of consistent scorings by the total
566 number of scorings, yielding an overall inter-rater reliability of 95%. Disagreements were
567 resolved by reaching consensus between the two coders.

568 **Analyses**

569 We first compared the DLD and TD groups' performance on the MSRT to examine
570 whether the two groups differed on this task. As the SR scores using the error method in syllable
571 for both groups and the SR scores using the binary method for the DLD group were not normally
572 distributed, Mann-Whitney U tests were conducted to compare the two groups' MSRT scores.
573 Receiver Operating Characteristic (ROC) curves were generated using SPSS v.26 to determine
574 the optimal cutoff point, sensitivity, and specificity for the MSRT. The ROC curve is a graph
575 which plots the sensitivity and specificity of a binary classification system as the discrimination
576 threshold (cutoff) varies (Fluss et al., 2005). Each ROC curve generates a value for the area
577 under the curve (AUC), which represents an overall estimate of the task's accuracy in classifying
578 individuals as with and without DLD. The AUC values are interpreted following the guidelines
579 in Swets et al. (2000): values between 0.90-1.0 are considered "excellent"; values between 0.80-
580 0.90 are considered "good", values between 0.70-0.80 are considered "fair", and values lower
581 than 0.70 are considered "poor". Following Redmond et al. (2019), the optimal cutoff points on

582 the ROC curves were identified using the Youden index (J) (Youden, 1950). The Youden index
583 J value captures the performance of a diagnostic test, and it is calculated following the formula
584 of $J = \text{sensitivity} + \text{specificity} - 1$. A J value of 0 indicates complete overlap between the affected
585 and unaffected groups and suggests that this classification task is useless. A J value of 1
586 indicates that the task could completely separate affected and unaffected groups. Therefore, the
587 optimal cutoff point should be associated with the maximum J value.

588 Once the optimal cutoff points and their associated sensitivity and specificity were
589 identified, the positive (LR+) and negative (LR-) likelihood ratios were calculated from
590 sensitivity and specificity values. LR+ and LR- respectively represents the probability that a
591 person with the condition testing positive for the condition and the probability that a person
592 without the condition testing negative for the condition. The likelihood ratios were calculated
593 using the following formula: $\text{LR+} = \text{sensitivity} / (1 - \text{specificity})$; $\text{LR-} = (1 - \text{sensitivity}) / \text{specificity}$.
594 We followed the guidelines specified in the introduction to interpret the sensitivity and
595 specificity values (Plante & Vance, 1994). Likelihood ratios were interpreted following
596 Dollaghan (2007): LR+ values between 3 and 10 indicate moderate positivity or a suggestive
597 level of clinical informativeness for identifying DLD; LR+ values at or above 10 are
598 confirmatory and clinically informative; LR- values between .1 and .2 indicate moderate
599 negativity; and LR- values at or below .1 are exclusionary and indicates high confidence in
600 ruling out a child with DLD.

601 As we intended for the MSRT to be used as a screening tool, greater emphasis is placed
602 on achieving high sensitivity (US Preventive Services Task Force, 2006). Children who perform
603 poorly on a screening tool would receive more comprehensive language assessments to confirm
604 their status. Therefore, this process would tolerate some false positives (inaccurately identify a

605 TD child as having DLD), as the misidentifications will be cleared up through further evaluation.
606 On the other hand, false negatives (inaccurately identify a child with DLD as TD) are not
607 desirable as they filter out children with DLD and stop them from receiving further evaluation
608 and intervention. A high sensitivity value indicates that there are few false negative results, thus
609 fewer cases of DLD are missed. Therefore, a good screening task would desire high sensitivity
610 values and may tolerate lower specificity values.

611 **Results**

612 *Group comparisons*

613 Figure 2 presents the two groups' SR scores. The middle line in each vertical box
614 represents the median, and the upper and lower lines represent the third and first quartile of the
615 data respectively. Using the error method in syllable, the DLD group received a mean score of
616 0.19 (SD=0.19) and the TD group received a mean score of 0.76 (SD=0.16). Using the binary
617 method, the DLD group received a mean score of 0.06 (SD=0.09) and the TD group received a
618 mean score of 0.62 (SD=0.21). Mann-Whitney U tests revealed that the TD group achieved
619 significantly higher scores compared to the DLD group using both scoring methods (error
620 method in syllable: $W=140$, $p<.001$, Cohen's $d=3.25$; binary method: $W=141$, $p<.001$, Cohen's
621 $d=3.47$).

622 *Classification accuracy*

623 The MSRT achieved AUC values of .984 and .982 using error method in syllable and
624 binary method, both demonstrating excellent classification accuracy. Using the error method in
625 syllable, an optimal cutoff score of .63 yielded sensitivity and specificity values of 100% and
626 87.5%. Using the binary method, an optimal cutoff score of .41 yielded the same sensitivity and

627 specificity values of 100% and 87.5%. The positive likelihood ratio is 8.0 using both scoring
628 methods, which demonstrate moderate classification of children with DLD. The negative
629 likelihood ratio is 0 using both scoring methods, which suggests high confidence in ruling out a
630 child with DLD.

631 **Discussion**

632 The current study aimed to design and validate a Mandarin sentence repetition task. In
633 study 1, we investigated the criterion validity of the MSRT by examining correlations between
634 TD children's performance on the MSRT and benchmark measures of language skills based on a
635 narrative task. Four narrative measures were derived, including two syntactic measures (MLU
636 and number of predicates), a vocabulary measure (VOCD), and a composite measure of four
637 linguistic structures (aspect markers, classifiers, passives, relative clauses). In addition, as
638 different scoring methods yielded distinct performance of SR tasks (Pham & Ebert, 2020;
639 Redmond, 2005; Stokes et al., 2006), we explored five potential scoring systems to determine the
640 best scoring method for the MSRT. In study 2, we calculated the sensitivity, specificity, and LRs
641 associated with the optimal cutoff score to examine whether the MSRT could accurately
642 differentiate between Mandarin-speaking preschoolers with and without DLD.

643 Study 1 showed that the MSRT is a valid measure to evaluate Mandarin-speaking
644 preschoolers' language ability. The two error methods, the binary method, and the percent of
645 correct syllable method significantly correlated with all validation measures, illustrating that the
646 MSRT in general could reflect children's language abilities as measured through narrative
647 sampling. However, scores using the percent of correct syllable method was right skewed and
648 resulted in a ceiling effect. This ceiling effect could be attributed to that the percent correct
649 syllable method did not consider the addition or the transposition of the syllables/words. In some

650 cases, a syllable corresponds to a stand-alone word in Mandarin, and word order is an important
651 consideration as children with DLD demonstrate difficulties with following the correct word
652 orders (Hansson & Nettelbladt, 1995). The core element scoring, which only considered the
653 accuracy of predefined grammatical targets, did not show significant correlations with MLU and
654 number of predicates. The core element scoring does not consider other elements in the
655 sentences except for the target structures, thus providing only a partial picture of children's
656 language abilities. In addition, the core element scoring resulted in a ceiling effect with over 50%
657 of children scoring 100% accurate on the task, lowering its ability to reveal individual
658 differences among children.

659 Study 2 showed that the MSRT can adequately differentiate Mandarin-speaking children
660 with and without DLD using both the error method in syllable and the binary scoring method.
661 Additional investigation in a new sample is needed to further compare the classification accuracy
662 associated with the two scoring methods. The AUCs derived from the ROC analysis
663 demonstrated excellent classification accuracy and the classification accuracy values all indicate
664 acceptable to good classification power. The 100% sensitivity using both scoring methods
665 further supported the MSRT as a language screening task for DLD. This study verifies the utility
666 of SR tasks to differentiate between children with and without DLD in a language that is
667 typologically distinct from the most studied Indo-European languages (e.g., Conti-Ramsden et
668 al., 2001; Leclercq et al., 2014; Redmond, 2005). In addition, our finding showed that the error
669 method in syllable and the binary method yielded the same classification accuracy in a Mandarin
670 SR task, which differs from the findings in other Asian languages wherein the error method
671 achieved higher overall classification accuracy compared to other scoring methods, including the
672 binary method (Cantonese: Stokes et al., 2006; Vietnamese: Pham & Ebert, 2020).

673 The high sensitivity of the MSRT could be attributed to our stimulus design that took into
674 consideration known areas of learning difficulties in Mandarin, including classifiers – noun-
675 modifying morphemes that are semantically complex (Hao et al., 2021), and two grammatical
676 features – passives (Zeng et al., 2018) and aspect markers (He & Sun, 2013). Errors on these
677 vulnerable structures were frequently observed in the responses produced by children with DLD.
678 For example, substitutions of specific classifiers with the general classifier *ge* (e.g., one *CL-kuai*
679 *chocolate* -> one *CL-ge chocolate*) and omissions of aspect markers (*ZAI* and *LE*) were quite
680 common. In addition, children with DLD changed passive sentences (e.g., *The wolf is defeated*
681 *by the smart goat.*) to either simple sentences (e.g., *The wolf defeated the smart goat.*) or BA-
682 sentences in Mandarin (e.g., *The wolf BA goat defeated.*). Future studies that more closely
683 examine the error patterns in children with DLD could further shed light on the linguistic
684 manifestations of DLD in Mandarin.

685 *Clinical Utility of the MSRT*

686 The design and validation of the MSRT remediates the paucity of Mandarin language
687 evaluation tools by providing clinicians and researchers with a quick screening tool to
688 differentiate children with and without DLD at an initial stage. As described earlier in the paper,
689 a high sensitivity value is desirable for a good screening task. The current MSRT is thus a
690 promising screening tool as the sensitivity is 100% and the specificity is 87.5% using both the
691 error (syllable) and binary methods. In addition, a screening task needs to be time-efficient so
692 that it can be given to a large number of individuals. For a vastly under-diagnosed disorder such
693 as DLD, this is especially important because the uncovering of the many hidden cases may
694 require universal screening. From our data collection experience, administering the MSRT takes
695 about six minutes and transcription and scoring using the error method take about eight minutes,

696 adding up to approximately 14 minutes to screen one child. When using the binary method, the
697 scoring can be completed online without transcriptions, which makes it around six minutes to
698 screen one child. The binary method may be selected when time is of the essence. The error
699 method in syllable may be selected when the examiner desires more in-depth information on the
700 child's error patterns. Moreover, the task does not require intensive training or considerable
701 linguistic expertise and can be administered and scored by classroom teachers or nurses. This
702 short and valid screening test will facilitate the allocation of limited resources to identify children
703 who are in need of further assessments and timely intervention.

704 It is important to use the optimal scoring methods and cutoff points in the real-life
705 application of MSRT as a screening tool for DLD. The sensitivity and specificity values
706 presented in this paper were associated with the reported optimal cutoff point (accuracy of 0.63
707 for error method, accuracy of 0.41 for binary method). Using a different cutoff may negatively
708 impact the classification accuracy of this task. In addition, we recruited the DLD group in a
709 hospital setting, and the need of medical assessment/consultation suggested parent concerns of
710 children's language and behavior. We recommend using a combination of the MSRT and a
711 measure of parent concern when viable, which may yield more accurate classification results
712 when using the MSRT as a universal screening tool.

713 **Limitations and future directions**

714 We employed a two-gate design instead of the more desired one-gate design in study 2.
715 In a one-gate design, all participants are recruited from a single population, whereas in a two-
716 gate design, affected and unaffected groups are recruited from separate populations. Although
717 only two of the 13 English diagnostic studies reviewed in Pawłowska (2014) used the one-gate
718 design in participant recruitment, the author emphasized the importance of a one-gate design in

719 diagnostic accuracy studies to avoid the influence from the fundamental differences across
720 populations. In study 2, we started with a one-gate design by recruiting both groups from a
721 hospital's outpatient clinic. However, this proved inefficient as we recruited only five TD
722 children during a period of nine months. We had to change course and recruited additional TD
723 children from the local community. Potential differences in sample characteristics may partly
724 contribute to the good classification accuracy of MSRT presented in this study.

725 We used an exploratory sample to examine the classification utility of the MSRT and the
726 corresponding cutoff score. Future studies should replicate the results with a confirmatory
727 sample of children with and without DLD to verify the classification power of this task and
728 further test the utility of the two scoring methods. A complete list of the sentences is provided in
729 appendix A. To achieve the one-gate design in participant recruitment, future studies could either
730 extend over a longer period in a hospital setting or carry out larger-scale screening in schools to
731 ensure that the participants are from the same population. Moreover, future studies could conduct
732 qualitative analysis of repetition response and compare children with and without DLD. Woon et
733 al. (2014) suggested that SR tasks are good candidates for qualitative evaluations and the errors
734 identified in SR tasks can help us understand the underlying deficits and particular weaknesses
735 associated with DLD in Mandarin-speaking populations.

736 **Acknowledgment**

737 The authors wish to thank the participating families for volunteering their time, the
738 research assistants at Nanjing Normal University and the clinician in the Developmental and
739 Behavioral Pediatrics Department in the Shanghai Children's Medical Center for administering
740 the tasks, Huanhuan Shi and Lue Shen for helping with the reliability check of the transcriptions

741 and coding, and Dr. Pumpki Lei Su for providing insightful comments to an earlier draft of the
742 paper.

743

744

745

References

746 Alloway, T. P., & Gathercole, S. (2005). Working memory and short-term sentence recall in
747 young children. *European Journal of Cognitive Psychology, 17*(2), 207-220.

748 Altman, C., Armon-Lotem, S., Fichman, S., & Walters, J. (2016). Macrostructure,
749 microstructure, and mental state terms in the narratives of English-Hebrew bilingual
750 preschool children with and without specific language impairment. *Applied*
751 *Psycholinguistics, 37*(1), 165.

752 American Psychiatric Association, & American Psychiatric Association. (2013). Diagnostic and
753 statistical manual of mental disorders: DSM-5. *Arlington, VA*.

754 Andreu, L., Sanz-Torrent, M., Guàrdia Olmos, J., & Macwhinney, B. (2011). Narrative
755 comprehension and production in children with SLI: An eye movement study. *Clinical*
756 *linguistics & phonetics, 25*(9), 767-783.

757 Archibald, L. M., & Joanisse, M. F. (2009). On the sensitivity and specificity of nonword
758 repetition and sentence recall to language and memory impairments in children. *Journal of*
759 *Speech, Language, and Hearing Research, 52*(4), 899-914.

- 760 Armon-Lotem, S., de Jong, J., & Meir, N. (Eds.). (2015). *Assessing multilingual children:*
761 *Disentangling bilingualism from language impairment*. Multilingual matters.
- 762 Armon-Lotem, S., & Meir, N. (2016). Diagnostic accuracy of repetition tasks for the
763 identification of specific language impairment (SLI) in bilingual children: evidence from
764 Russian and Hebrew. *International journal of language & communication disorders*, 51(6),
765 715-731.
- 766 Baddeley, A. (2000). Working memory and language processing. *Benjamins translation*
767 *library*, 40, 1-16.
- 768 Bedore, L. M., & Leonard, L. B. (2001). Grammatical morphology deficits in Spanish-speaking
769 children with specific language impairment. *Journal of Speech, Language, and Hearing*
770 *Research*.
- 771 Bishop, D. V. (1997). Cognitive neuropsychology and developmental disorders: Uncomfortable
772 bedfellows. *The Quarterly Journal of Experimental Psychology Section A*, 50(4), 899-923.
- 773 Bishop, D. V., Snowling, M. J., Thompson, P. A., Greenhalgh, T., & Catalise Consortium.
774 (2016). CATALISE: A multinational and multidisciplinary Delphi consensus study.
775 Identifying language impairments in children. *PLOS one*, 11(7), e0158753.
- 776 Bishop, D. V., Snowling, M. J., Thompson, P. A., Greenhalgh, T., Catalise-2 Consortium,
777 Adams, C., ... & Boyle, C. (2017). Phase 2 of CATALISE: A multinational and
778 multidisciplinary Delphi consensus study of problems with language development:
779 Terminology. *Journal of Child Psychology and Psychiatry*, 58(10), 1068-1080.

- 780 Blom, E., Boerma, T., & de Jong, J. (2020). Multilingual Assessment Instrument for Narratives
781 (MAIN) adapted for use in Dutch. *ZAS Papers in Linguistics*, 64, 51-56.
- 782 Bortolini, U., Caselli, M. C., & Leonard, L. B. (1997). Grammatical deficits in Italian-speaking
783 children with specific language impairment. *Journal of Speech, Language, and Hearing*
784 *Research*, 40(4), 809-820.
- 785 Boudreau, D. (2008). Narrative abilities: Advances in research and implications for clinical
786 practice. *Topics in Language Disorders*, 28(2), 99-114.
- 787 Chan, A., Cheng, K., Kan, R., Wong, A. M. Y., Fung, R., Wong, J., ... & Gagarina, N. (2020).
788 The Multilingual Assessment Instrument for Narratives (MAIN): Adding Cantonese to
789 MAIN. *ZAS Papers in Linguistics*, 64, 23-29.
- 790 Chao, Y. R., & Yang, L. S. (1949). *Concise Dictionary of Spoken Chinese*, Cambridge, Mass.:
791 Harvard University Press.
- 792 Chen, J. H., & Chen, H. Y. (2000). Manual for the Wechsler preschool and primary scale of
793 intelligence-revised. *Taipei, Taiwan: Chinese Behavioral Science Corporation*.
- 794 Chéng Y. (2003). *汉语字基语法: 语素层造句的理论和实践* [A zì-based grammar of Chinese:
795 theory and practice in building sentences with morphemes], Shanghai: Fudan University
796 Press.
- 797 Christensen, R. V. (2019). Sentence Repetition: A Clinical Marker for Developmental Language
798 Disorder in Danish. *Journal of Speech, Language and Hearing Research (Online)*, 62(12),
799 4450-4463.

- 800 Clarke, M. G., & Leonard, L. B. (1996). Lexical comprehension and grammatical deficits in
801 children with specific language impairment. *Journal of communication disorders*, 29(2),
802 95-105.
- 803 Conti-Ramsden, G., Botting, N., & Faragher, B. (2001). Psycholinguistic markers for specific
804 language impairment (SLI). *Journal of child psychology and psychiatry*, 42(6), 741-748.
- 805 Devescovi, A., & Cristina Caselli, M. (2007). Sentence repetition as a measure of early
806 grammatical development in Italian. *International Journal of Language & Communication*
807 *Disorders*, 42(2), 187-208.
- 808 Dietz, A., & Boyle, M. (2018). Discourse measurement in aphasia research: have we reached the
809 tipping point?. *Aphasiology*, 32(4), 459-464.
- 810 Dollaghan, C. A. (2007). *The handbook for evidence-based practice in communication disorders*.
811 Paul H Brookes Publishing.
- 812 Dong, J. (2010). The enregisterment of Putonghua in practice. *Language &*
813 *Communication*, 30(4), 265-275.
- 814 Duanmu, S., & Dong, Y. (2016). Elastic words in Chinese. In *The Routledge Encyclopedia of the*
815 *Chinese Language* (pp. 490-506). Routledge.
- 816 Ehrlert, D. J., & McGhee, R. L. (2008). *PTONI: Primary test of nonverbal intelligence*. Austin,
817 TX: Pro-Ed.
- 818 Eisenberg, S. L. (2020). Using general language performance measures to assess grammar
819 learning. *Topics in Language Disorders*, 40(2), 135-148.

- 820 Finestack, L. H., & Satterlund, K. E. (2018). Current practice of child grammar intervention: A
821 survey of speech-language pathologists. *American Journal of Speech-Language*
822 *Pathology*, 27(4), 1329-1351.
- 823 Fleckstein, A., Prévost, P., Tuller, L., Sizaret, E., & Zebib, R. (2018). How to identify SLI in
824 bilingual children: a study on sentence repetition in French. *Language Acquisition*, 25(1),
825 85-101.
- 826 Fletcher, P., Leonard, L. B., Stokes, S. F., & Wong, A. M. Y. (2005). The expression of aspect in
827 Cantonese-speaking children with specific language impairment. *Journal of Speech,*
828 *Language, and Hearing Research*.
- 829 Fluss, R., Faraggi, D., & Reiser, B. (2005). Estimation of the Youden Index and its associated
830 cutoff point. *Biometrical Journal: Journal of Mathematical Methods in Biosciences*, 47(4),
831 458-472.
- 832 Gagarina, N. V., Klop, D., Kunnari, S., Tantele, K., Välimaa, T., Balčiūnienė, I., ... & Walters,
833 J. (2012). MAIN: Multilingual assessment instrument for narratives. *ZAS papers in*
834 *linguistics*, 56, 155-155.
- 835 Gordon Jr, R. G. (2005). Ethnologue, languages of the world. [http://www. Ethnologue. Com/](http://www.ethnologue.com/)
- 836 Hansson, K., & Nettelbladt, U. (1995). Grammatical characteristics of Swedish children with
837 SLI. *Journal of Speech, Language, and Hearing Research*, 38(3), 589-598.
- 838 Hao, Y., Sheng, L., Zhang, Y., Jiang, F., de Villiers, J., Lee, W., & Liu, X. L. (2018). A narrative
839 evaluation of Mandarin-speaking children with language impairment. *Journal of Speech,*
840 *Language, and Hearing Research*, 61(2), 345-359.

- 841 He, X., Sun, L. (2013). 汉语特殊性语言障碍儿童体标记“了”和“在”的产出研究 [The
842 production of aspect markers “le” and “zai” by Mandarin-speaking children with specific
843 language impairment]. *Journal of Foreign Language Education*, 34(2), 27-32.
- 844 He, W., Xu, N., & Ji, R. (2017). Effects of age and location in Chinese relative clauses
845 processing. *Journal of psycholinguistic research*, 46(5), 1067-1086.
- 846 Hwang, M. (2012). Sentence repetition as a clinical marker of specific language impairment in
847 Korean-speaking preschool children. *Communication Sciences & Disorders*, 17(1), 1-14.
- 848 Jessup, B., Ward, E., Cahill, L., & Keating, D. (2008). Teacher identification of speech and
849 language impairment in kindergarten students using the Kindergarten Development
850 Check. *International journal of speech-language pathology*, 10(6), 449-459.
- 851 Jespersen, O. (1922). *Language: Its nature, development and origin* (Vol. 68). H. Holt.
- 852 Joanisse, M. F., & Seidenberg, M. S. (1998). Specific language impairment: A deficit in
853 grammar or processing?. *Trends in cognitive sciences*, 2(7), 240-247.
- 854 Johnston, J. R., & Kamhi, A. G. (1984). Syntactic and semantic aspects of the utterances of
855 language-impaired children: The same can be less. *Merrill-Palmer Quarterly (1982-)*, 65-
856 85.
- 857 Justice, L. M., Bowles, R., Pence, K., & Gosse, C. (2010). A scalable tool for assessing
858 children’s language abilities within a narrative context: The NAP (Narrative Assessment
859 Protocol). *Early Childhood Research Quarterly*, 25(2), 218-234.
- 860 Kan, P. F., & Windsor, J. (2010). Word learning in children with primary language impairment:
861 A meta-analysis.

- 862 Kennedy, G. A. (1951). The monosyllabic myth. *Journal of the American Oriental*
863 *Society*, 71(3), 161-166.
- 864 Klem, M., Melby-Lervåg, M., Hagtvet, B., Lyster, S. A. H., Gustafsson, J. E., & Hulme, C.
865 (2015). Sentence repetition is a measure of children's language skills rather than working
866 memory limitations. *Developmental science*, 18(1), 146-154.
- 867 Knox, E., & Conti-Ramsden, G. (2003). Bullying risks of 11-year-old children with specific
868 language impairment (SLI): Does school placement matter?. *International Journal of*
869 *Language & Communication Disorders*, 38(1), 1-12.
- 870 Komeili, M., & Marshall, C. R. (2013). Sentence repetition as a measure of morphosyntax in
871 monolingual and bilingual children. *Clinical linguistics & phonetics*, 27(2), 152-162.
- 872 Laws, G., Briscoe, J., Ang, S. Y., Brown, H., Hermena, E., & Kapikian, A. (2015). Receptive
873 vocabulary and semantic knowledge in children with SLI and children with Down
874 syndrome. *Child Neuropsychology*, 21(4), 490-508.
- 875 Leclercq, A. L., Quémart, P., Magis, D., & Maillart, C. (2014). The sentence repetition task: A
876 powerful diagnostic tool for French children with specific language impairment. *Research*
877 *in developmental disabilities*, 35(12), 3423-3430.
- 878 Leonard, L. B., Miller, C., & Gerber, E. (1999). Grammatical morphology and the lexicon in
879 children with specific language impairment. *Journal of Speech, Language, and Hearing*
880 *Research*, 42(3), 678-689.
- 881 Leonard, L. B., Kas, B., & Pléh, C. (2009). The use of tense and agreement by Hungarian-
882 speaking children with language impairment.

- 883 Leonard, L. B. (2014). Specific language impairment across languages. *Child development*
884 *perspectives*, 8(1), 1-5.
- 885 Levorato, C., & Roch, M. (2020). Italian adaptation of the Multilingual Assessment Instrument
886 for Narratives. *ZAS Papers in Linguistics*, 64, 139-146.
- 887 Li, C. N., & Thompson, S. A. (1989). *Mandarin Chinese: A functional reference grammar* (Vol.
888 3). Univ of California Press.
- 889 Lín, H. (1952). 汉语是不是单音节语? [Is Chinese a monosyllabic language?], *中国语文* 11,
890 6-11.
- 891 Liu, X. L., de Villiers, J., Ning, C., Rolfhus, E., Hutchings, T., Lee, W., ... & Zhang, Y. W.
892 (2017). Research to establish the validity, reliability, and clinical utility of a comprehensive
893 language assessment of Mandarin. *Journal of Speech, Language, and Hearing*
894 *Research*, 60(3), 592-606.
- 895 Luo, J., Yang, W., Chan, A., Cheng, K., Kan, R., & Gagarina, N. (2020). The Multilingual
896 Assessment Instrument for Narratives (MAIN): Adding Mandarin to MAIN. *ZAS Papers in*
897 *Linguistics*, 64, 159-162.
- 898 McGregor, K. K., Newman, R. M., Reilly, R. M., & Capone, N. C. (2002). Semantic
899 representation and naming in children with specific language impairment.
- 900 MacWhinney, B. (2000). *The CHILDES Project: Tools for analyzing talk*. Psychology Press.
- 901 Marinis, T., & Armon-Lotem, S. (2015). Sentence repetition. *Assessing multilingual children:*
902 *Disentangling bilingualism from language impairment*, 95-124.

- 903 McArthur, G. M., Hogben, J. H., Edwards, V. T., Heath, S. M., & Mengler, E. D. (2000). On the
904 “specifics” of specific reading disability and specific language impairment. *Journal of*
905 *Child Psychology and Psychiatry*, 41(7), 869-874.
- 906 Meir, N., Walters, J., & Armon-Lotem, S. (2016). Disentangling SLI and bilingualism using
907 sentence repetition tasks: The impact of L1 and L2 properties. *International Journal of*
908 *Bilingualism*, 20(4), 421-452.
- 909 Newcomer, P., & Hammill, D. (2019). *Told-p: 5: Test of Language Development. Primary.*
910 Austin, TX: Pro-Ed.
- 911 Ning, C. Y., Liu, X. L., & de Villiers, J. G. (2014). The diagnostic receptive and expressive
912 assessment of Mandarin. *Dallas, TX: Bethel Hearing and Speaking Training Center.*
- 913 Norbury, C. F., Gooch, D., Wray, C., Baird, G., Charman, T., Simonoff, E., ... & Pickles, A.
914 (2016). The impact of nonverbal ability on prevalence and clinical presentation of language
915 disorder: evidence from a population study. *Journal of Child Psychology and*
916 *Psychiatry*, 57(11), 1247-1257.
- 917 Paul, R., Norbury, C.F., & Gosse, C. (2017). *Language disorders from infancy through*
918 *adolescence: Listening, speaking, reading, writing and communication* (5th ed.). St. Louis,
919 MO: Mosby Elsevier.
- 920 Pawłowska, M. (2014). Evaluation of three proposed markers for language impairment in
921 English: A meta-analysis of diagnostic accuracy studies. *Journal of Speech, Language, and*
922 *Hearing Research*, 57(6), 2261-2273.

- 923 Pham, G., & Ebert, K. D. (2020). Diagnostic Accuracy of Sentence Repetition and Nonword
924 Repetition for Developmental Language Disorder in Vietnamese. *Journal of Speech,
925 Language, and Hearing Research, 63*(5), 1521-1536.
- 926 Plante, E., & Vance, R. (1994). Selection of preschool language tests: A data-based
927 approach. *Language, Speech, and Hearing Services in Schools, 25*(1), 15-24.
- 928 Polišenská, K., Chiat, S., & Roy, P. (2015). Sentence repetition: What does the task measure?.
929 *International Journal of Language & Communication Disorders, 50*(1), 106-118
- 930 Redmond, S. M. (2005). Differentiating SLI from ADHD using children's sentence recall and
931 production of past tense morphology. *Clinical Linguistics & Phonetics, 19*(2), 109-127.
- 932 Redmond, S. M., Ash, A. C., Christopoulos, T. T., & Pfaff, T. (2019). Diagnostic accuracy of
933 sentence recall and past tense measures for identifying children's language
934 impairments. *Journal of Speech, Language, and Hearing Research, 62*(7), 2438-2454.
- 935 Rezzonico, S., Chen, X., Cleave, P. L., Greenberg, J., Hipfner-Boucher, K., Johnson, C. J., ... &
936 Girolametto, L. (2015). Oral narratives in monolingual and bilingual preschoolers with
937 SLI. *International Journal of Language & Communication Disorders, 50*(6), 830-841.
- 938 Rice, M. L. (2013). Language growth and genetics of specific language
939 impairment. *International Journal of Speech-Language Pathology, 15*(3), 223-233.
- 940 Richards, B. (1987). Type/token ratios: What do they really tell us?. *Journal of child
941 language, 14*(2), 201-209.
- 942 Richards, B., & Malvern, D. (2000). Measuring vocabulary richness in teenage learners of
943 French.

- 944 Rispens, J. E. (2004). *Syntactic and phonological processing in developmental dyslexia*.
945 Groningen: Rijksuniversiteit Groningen, Faculteit der Letteren.
- 946 Semel, E., Wiig, E., & Secord, W. (1994). *Clinical Evaluation of Language Fundamentals-*
947 *Revised*. San Antonio, TX: The Psychological Corporation.
- 948 Sheng L. (2014). Semantic Development in Children with Language Impairments. *Encyclopedia*
949 *of language development*, 534-538.
- 950 Sheng, L., & McGregor, K. K. (2010). Lexical–semantic organization in children with specific
951 language impairment. *Journal of Speech, Language, and Hearing Research*, 53(1), 146-
952 159.
- 953 Sheng, L., Shi, H., Wang, D., Hao, Y., & Zheng, L. (2020). Narrative Production in Mandarin-
954 Speaking Children: Effects of Language Ability and Elicitation Method. *Journal of Speech,*
955 *Language, and Hearing Research*, 63(3), 774-792.
- 956 Sheng, L., Su, P.L., Wang, D., Yu, J., Lu, T.-H., Shen, L., Hao, Y., & Lam, B.P.W. (in
957 preparation). Manifestations of developmental language disorder in Chinese children: A
958 systematic review and meta-analysis.
- 959 Spencer, E., Bryant, L., & Colyvas, K. (2020). Minimizing variability in language sampling
960 analysis: A practical way to calculate text length and time variability and measure reliable
961 change when assessing clients. *Topics in Language Disorders*, 40(2), 166-181.
- 962 Stokes, S. F., Wong, A. M., Fletcher, P., & Leonard, L. B. (2006). Nonword repetition and
963 sentence repetition as clinical markers of specific language impairment: The case of
964 Cantonese. *Journal of Speech, Language, and Hearing Research*, 49(2), 219-236.

- 965 Sung, Y. T., Cha, J. H., Tu, J. Y., Wu, M. D., & Lin, W. C. (2016). Investigating the processing
966 of relative clauses in Mandarin Chinese: evidence from eye-movement data. *Journal of*
967 *psycholinguistic research*, 45(5), 1089-1113.
- 968 Swets, J. A., Dawes, R. M., & Monahan, J. (2000). Psychological science can improve
969 diagnostic decisions. *Psychological science in the public interest*, 1(1), 1-26.
- 970 Tardif, T., & Fletcher, P. (2008). Chinese Communicative Development Inventories: user's
971 guide and manual.
- 972 Thal, D. J., O'Hanlon, L., Clemmons, M., & Fralin, L. (1999). Validity of a parent report
973 measure of vocabulary and syntax for preschool children with language
974 impairment. *Journal of Speech, Language, and Hearing Research*, 42(2), 482-496.
- 975 Thompson, S. A., & Tao, H. (2010). Conversation, grammar, and fixedness: adjectives in
976 Mandarin revisited. *Chinese Language and Discourse*, 1(1), 3-30.
- 977 Tomblin, J. B., Records, N. L., & Zhang, X. (1996). A system for the diagnosis of specific
978 language impairment in kindergarten children. *Journal of Speech, Language, and Hearing*
979 *Research*, 39(6), 1284-1294.
- 980 Tomblin, J. B., Records, N. L., Buckwalter, P., Zhang, X., Smith, E., & O'Brien, M. (1997).
981 Prevalence of specific language impairment in kindergarten children. *Journal of speech,*
982 *language, and hearing research*, 40(6), 1245-1260.
- 983 Trauner, D., Wulfeck, B., Tallal, P., & Hesselink, J. (2000). Neurological and MRI profiles of
984 children with developmental language impairment. *Developmental Medicine & Child*
985 *Neurology*, 42(7), 470-475.

- 986 Tsimpli, I. M., Andreou, M., & Peristeri, E. (2020). The multilingual assessment instrument for
987 narratives: Greek. *ZAS Papers in Linguistics*, 64, 101-106.
- 988 Tweedie, F. J., & Baayen, R. H. (1998). How variable may a constant be? Measures of lexical
989 richness in perspective. *Computers and the Humanities*, 32(5), 323-352.
- 990 US Preventive Services Task Force. (2006). Screening for speech and language delay in
991 preschool children: recommendation statement. *Pediatrics*, 117(2), 497-501.
- 992 Van der Lely, H. K., & Battell, J. (2003). Wh-movement in children with grammatical SLI: A
993 test of the RDDR hypothesis. *Language*, 153-181.
- 994 Vân Hoàng, T., Schelstraete, M. A., Trần, Q. D., & Bragard, A. (2014). La répétition de
995 phrases en vietnamien—un marqueur des troubles du langage oral et des troubles du
996 comportement Sentence repetition in Vietnamese—a marker of oral language and behavioral
997 difficulties. *Canadian Journal of Speech-Language Pathology and Audiology*, 37(4), 280-
998 297.
- 999 Wang, W. S., & Sun, C. (2015). *The Oxford handbook of Chinese linguistics*. Oxford University
1000 Press.
- 1001 Watkins, R. V., Kelly, D. J., Harbers, H. M., & Hollis, W. (1995). Measuring children's lexical
1002 diversity: Differentiating typical and impaired language learners. *Journal of Speech,*
1003 *Language, and Hearing Research*, 38(6), 1349-1355.
- 1004 Wechsler, D. (1989). Wechsler preschool and primary intelligence scales for children,
1005 revised. New York, NY, The Psychological Corporation.
- 1006 Wiedenhof, J. (2015). *A grammar of Mandarin*. John Benjamins Publishing Company.

- 1007 Wiig, E., Semel, E., & Seccord, W. (2013). Clinical Evaluation of Language Fundamentals-
1008 Revised. *NY: Merrill Divison of Macmillan.*
- 1009 Woon, C. P., Yap, N. T., Lim, H. W., & Wong, B. E. (2014). Measuring Grammatical
1010 Development in Bilingual Mandarin-English Speaking Children with a Sentence Repetition
1011 Task. *Journal of Education and Learning, 3(3)*, 144-157.
- 1012 Youden, W. J. (1950). Index for rating diagnostic tests. *Cancer, 3(1)*, 32-35.
- 1013 Zeng, T., Zhu, T., Li, X., & Zhu, R. (2018). Passive Structure Features in Mandarin-Speaking
1014 Children with Specific Language Impairment: Optional Movement. *Journal of Language,
1015 Linguistics, and Literature, 4(1)*, 8-18.
- 1016
- 1017
- 1018
- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025
- 1026

1027

1028

1029

Tables and Figures

1030 Table 1. Correlations between SR scoring systems and validation measures

	MLU	VOCD	Number of predicates	Structure Composite
Binary	.435*	.356*	.380*	.594*
Core element	.246	.342*	.148	.397*
Error method – syllable	.432*	.402*	.389*	.587*
Error method – word	.419*	.400*	.392*	.599*
Percent of correct syllable	.383*	.441*	.379*	.559*

1031 Note: p level is corrected using the Bonferroni correction, and the corrected p value is .003

1032 * p<.003

1033 Table 2. The demographic characteristics and standardized test scores of the TD and DLD
 1034 groups

1035

Measure	TD (N=16)		DLD (N=16)		t/W	p value	Cohen's d
	Mean (SD)	Range	Mean (SD)	Range			
Age	60.8 (6.5)	50-71	60.3 (6.8)	50-70	$t=.10$.92	.05
Mandarin Exposure	3.8 (1.0)	1-5	3.9 (1.1)	1-5	$W=254$.68	.06
Maternal Education	4.2 (.8)	3-5	3.8 (.9)	3-5	$W=232$.19	.52
Performance IQ	120.9 (11.7)	99-139	101.9 (11.6)	87-124	$t=4.5$	<.001	1.63
DREAM total score	109.7 (8.2)	94-123	84.8 (7.6)	72-103	$t=8.9$	<.001	3.14

1036 Note: Value t is reported for t-tests when the two variables under comparison were normally
 1037 distributed; W is reported for Mann Whitney U tests when one or more of the variables under
 1038 comparison were not normally distributed. Age is reported in months. Mandarin exposure is
 1039 reported on a scale from 1 to 5: 1 means <19%; 2 means 20%-39%; 3 means 40%-59%, 4 means
 1040 60%-79%, and 5 means 80%-100%. One TD child and one child with DLD received a Mandarin
 1041 exposure score of 1. The TD child had exposure to Shanghai dialect at home and attended
 1042 English classes four to five hours each week. The child with DLD had exposure to both Shanghai
 1043 and Henan dialects at home and attended English classes for half a year. Performance IQ is
 1044 measured by WPPSI-R and is reported as standard scores. Maternal education is reported on a
 1045 scale from 1 to 5 through a parent questionnaire: 1 means middle school or lower, 2 means high
 1046 school, 3 means some college, 4 means bachelor's degree, and 5 means master's degree. The
 1047 total scores on DREAM are reported as standard scores.

1048

1049

1050

1051

1052

1053

1054

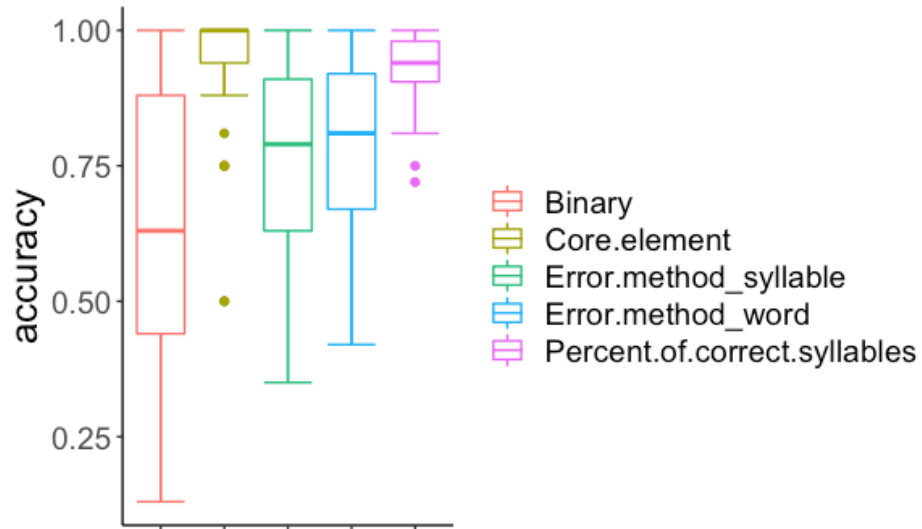
1055

1056

1057 Figure captions:

1058 Figure 1. Children's SR scores using the five different scoring methods

1059 Figure 2. Children's MSRT accuracy score by group and scoring method

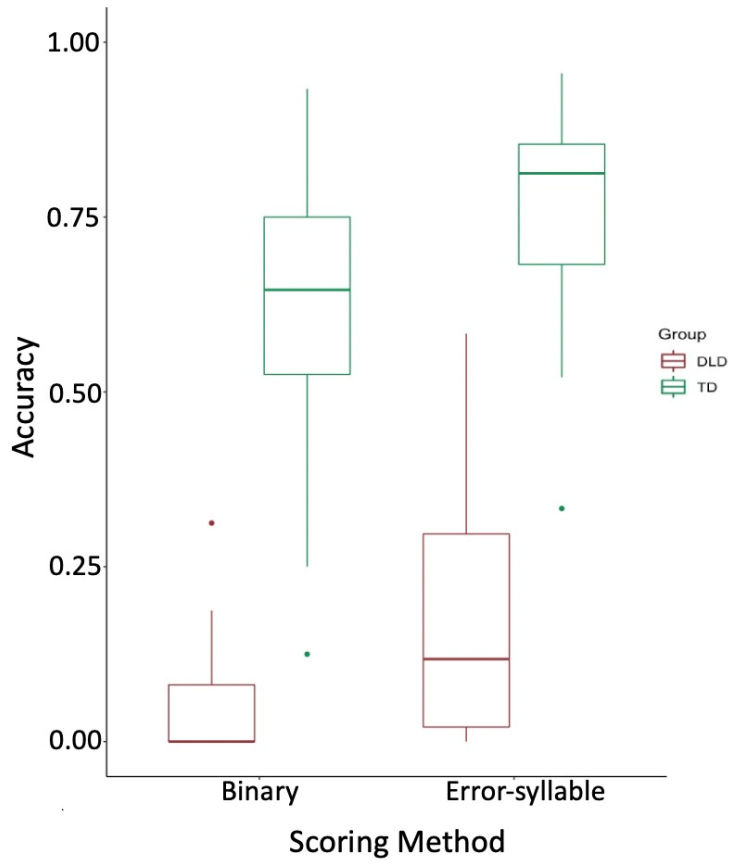


26 Figure 1. Children's SR scores using the five different scoring methods

27

28

29



1060 Appendix A. Sentences in the MSRT

1061 Classifiers are bolded and relative clauses are underlined.

1062 PASS=passive marker; SFP=sentence final particle; ProM = progressive marker; PerM =
1063 perfective marker; LP=linking particle; RP=resultative particle。

1064 Passive Sentences:

- 1065 1. 那 只 白色 的 小狗 被 妈妈 抱 走 了。
1066 That CL white LP dog PASS mom carry away (RP) SFP
1067 (That white dog was carried away by mom.)
- 1068 2. 那 件 破旧 的 毛衣 被 姐姐 扔 掉 了。
1069 That CL old LP sweater PASS sister throw away (RP) SFP
1070 (That old sweater was thrown away by sister.)
- 1071 3. 男孩 被 那 个 穿 裙子 的 女孩 绊 倒 了。
1072 Boy PASS that CL wear dress LP girl tripped fell (RP) SFP
1073 (The boy was tripped by that girl who wore a dress.)
- 1074 4. 大灰狼 被 那 只 聪明 的 小羊 打 败 了。
1075 Wolf PASS that CL smart LP sheep defeated lost (RP) SFP
1076 (The wolf was defeated by that smart sheep.)
- 1077 5. 蛋糕 被 那 个 戴 眼镜 的 女孩 吃 光 了。
1078 Cake PASS that CL wear glasses LP girl eaten up (RP) SFP
1079 (The cake was eaten by the girl who wore glasses.)
- 1080 6. 那 辆 红色 的 自行车 被 小偷偷 走 了。
1081 That CL red LP bicycle PASS thief stole away (RP) SFP
1082 (That red bicycle was stolen by the thief.)
- 1083 7. 那 个 长头发 的 女孩 被 一 只 猫 抓 伤 了。
1084 That CL long hair LP girl PASS one CL cat scratched hurt (RP) SFP
1085 (That girl with long hair was scratched by a cat.)
- 1086 8. 棒棒糖 被 那 个 高个子 的 男孩 抢 走 了。
1087 Lollipop PASS that CL tall LP boy robbed away (RP) SFP
1088 (The lollipop was robbed by that tall boy.)

1089 Aspect marker sentences:

- 1090 1. 那 个 戴 帽子 的 男孩 在 骑 自行车。
1091 That CL wear hat LP boy ProM ride bicycle
1092 That boy who wears a hat is riding a bicycle.
- 1093 2. 哥哥 在 组装 那 辆 绿色 的 玩具车。
1094 Older brother ProM assemble that CL green LP toy car
1095 Older brother is assembling that green toy car.
- 1096 3. 那 群 中班 的 小朋友 在 高兴 地 荡秋千。
1097 Those CL kindergarten LP children ProM happily LP play on a swing
1098 Those kindergarten children are playing on a swing happily.
- 1099 4. 老师 在 严厉 地 批评 那 个 淘气 的 男孩。
1100 Teacher ProM harshly LP criticize that CL naughty LP boy

1101 The teacher is harshly criticizing that naughty boy.

1102 LE:

1103 1. 小 老鼠 吃 了 一 块 美 味 的 巧 克 力。

1104 Little mouse eat PerM one CL delicious LP chocolate

1105 The little mouse has eaten a piece of delicious chocolate.

1106 2. 那 个 背 书 包 的 男 孩 掉 了 一 本 书。

1107 That CL carry bag LP boy drop PerM one CL book

1108 That boy who carries a bag has dropped a book.

1109 3. 那 个 可 爱 的 女 孩 大 声 地 唱 了 一 首 歌。

1110 That CL cute LP girl loudly LP sing PerM one CL song

1111 That cute girl has sung a song loudly.

1112 4. 妈 妈 认 真 地 洗 了 一 件 漂 亮 的 衣 服。

1113 Mother seriously LP wash PerM one CL beautiful LP clothes

1114 Mother has seriously washed a beautiful clothes.

1115

1116

1117

A. Scatter plots between the SR scores and the narrative measures in study 1

Figure 1. Scatter Plots of the Error Method Word Scoring and Narrative Measures

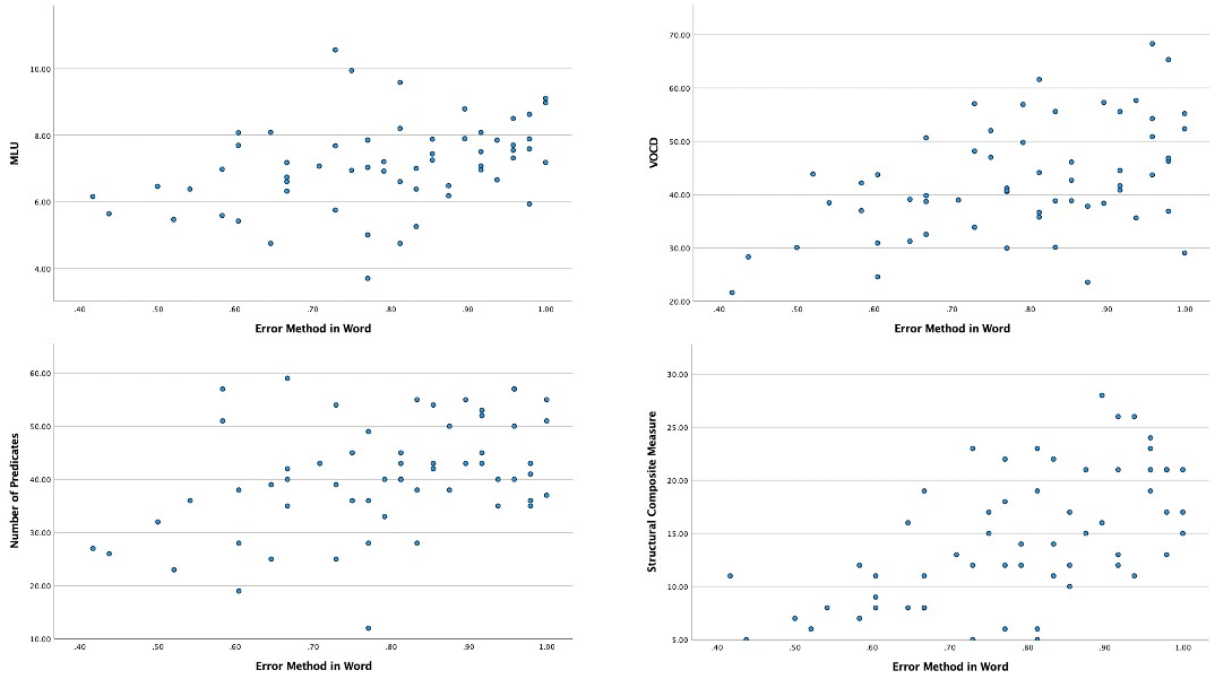


Figure 2. Scatter Plots of the Error Method Syllable Scoring and Narrative Measures

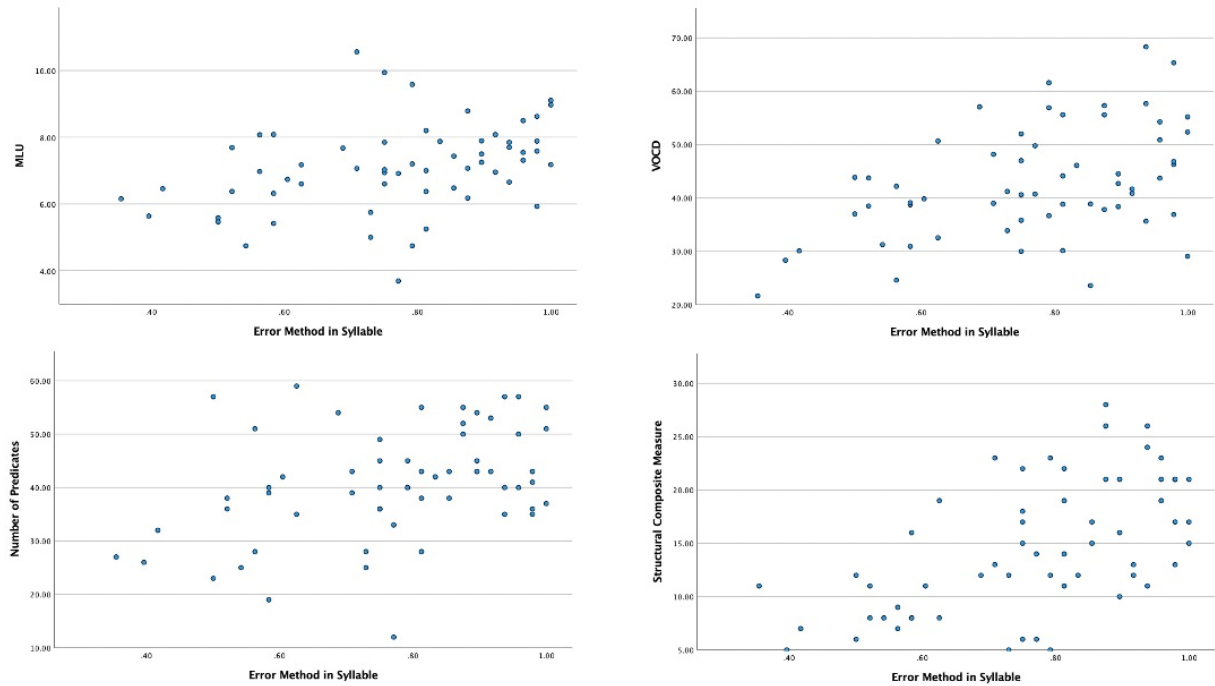


Figure 3. Scatter Plots of the Binary Scoring and Narrative Measures

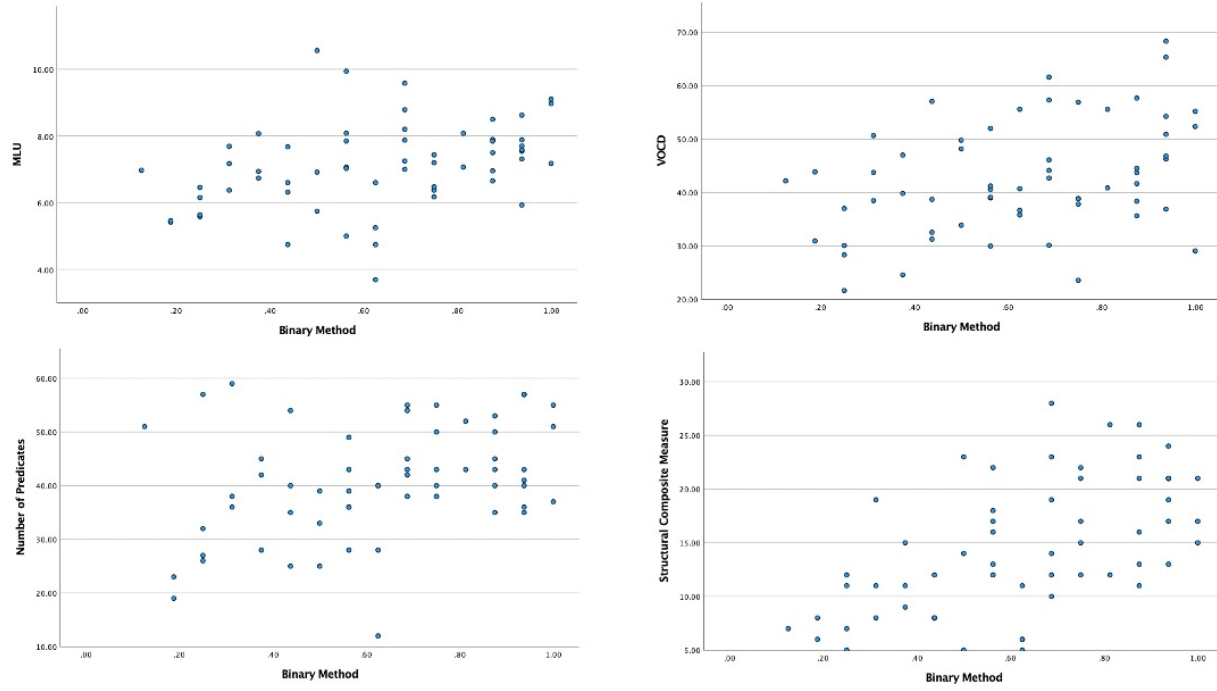


Figure 4. Scatter Plots of the Core Element Scoring and Narrative Measures

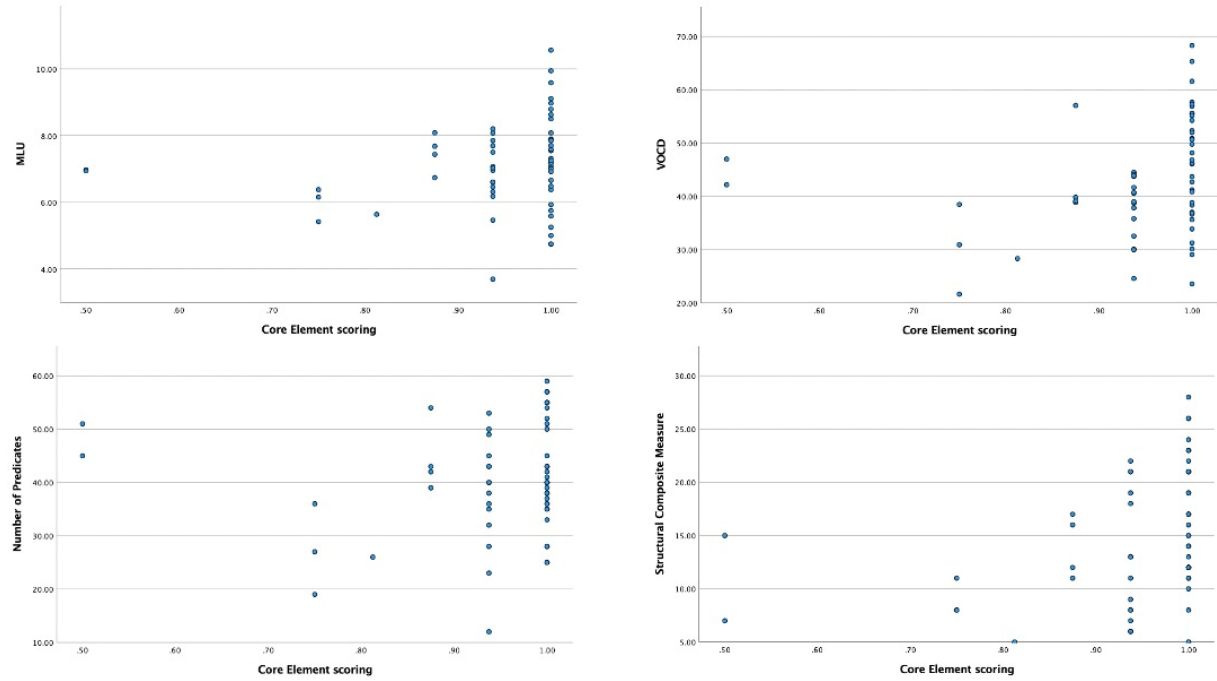
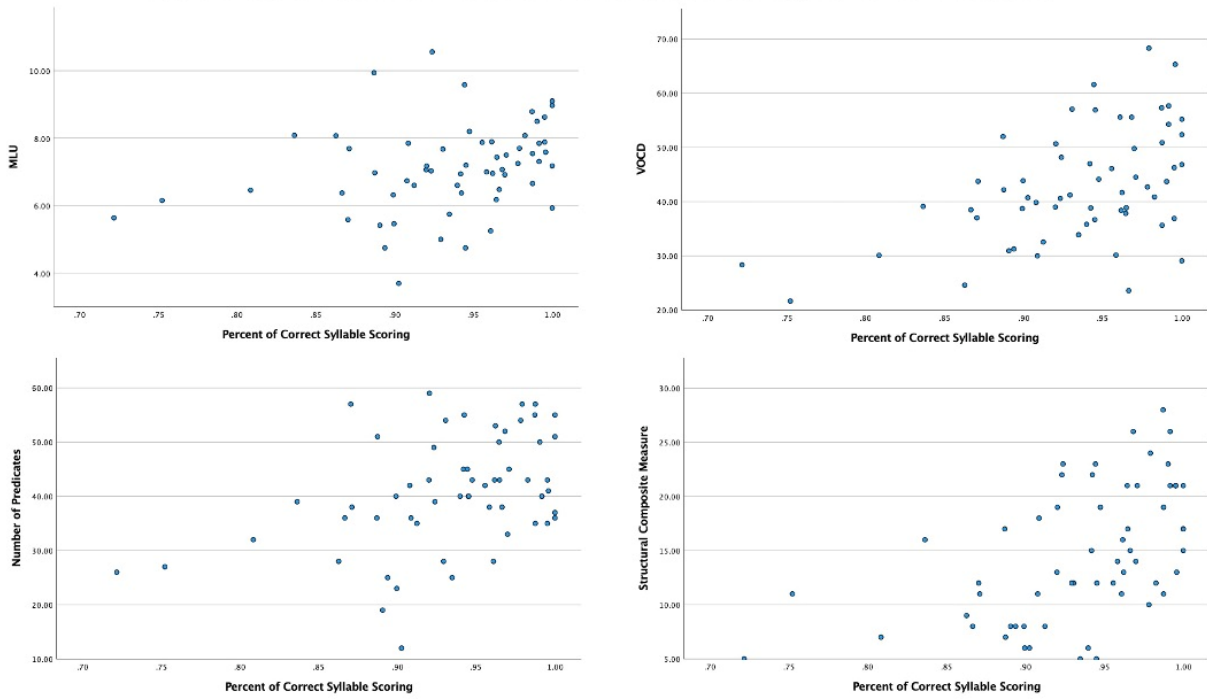


Figure 5. Scatter Plots of the Percent of Correct Syllable Scoring and Narrative Measures



B. Examples of a sample response and its scoring procedures

Target sentence:

小 老鼠 吃 了 一 块 美味的 巧克力。
xiao3 lao3shu3 chi1 le yi1 kuai4 mei3wei4de qiao3ke1li4
Little mouse eat PerM one CL delicious chocolate
The little mouse has eaten a piece of delicious chocolate.

Child's response:

小 老鼠 吃 了 一 块 好吃的 巧克力。
xiao3 lao3shu3 chi1 le yi1 kuai4 hao3chi1de qiao3ke1li4
Little mouse eat PerM one CL *tasty* chocolate
The little mouse has eaten a piece of *tasty* chocolate.

Error method – word: The child substituted the word “美味” with “好吃”, which was considered as one word error. A score of **2** (1 error) was given for this response using the error method in word system.

Error method – syllable: The child substituted two syllables “美味” with “好吃”, which was considered as two syllable errors. A score of **1** (two to three errors) was given for this response using the error method in syllable system.

Binary method: As this response was not a completely accurate repetition of the target, a score of **0** was given using the binary method.

Core element scoring: The core element in the target sentence is the aspect marker structure (verb + aspect marker + noun). The target core element is present in the child's response, so a score of **1** was given using the core element scoring.

Percent of correct syllable scoring: The target sentence has a total of 13 syllables and 11 of them were repeated correctly in the response. The percent of correct syllable scoring thus gives a score of 11/13, which is **0.85**.

C. Descriptive Data of Children's Scores on MSRT and Narrative Measures in Study 1

Measure	Range	Mean	SD
Error method – word	0.42-1.00	0.79	0.15
Error method – syllable	0.35-1.00	0.77	0.17
Binary method	0.13-1.00	0.63	0.25
Core element method	0.50-1.00	0.94	0.11
Percent of correct syllable method	0.72-1.00	0.93	0.06
MLU-w	3.69-10.56	7.11	1.30
VOCD	21.65-68.32	42.75	10.39
Number of predicates	12-59	40.69	10.40
Structural composite	5-28	14.63	6.14

D. Examples of the narrative measures

Structure type	Subtype	Example
Predicate	Verb	然后 小 鸟 妈妈 爬 到 树上。 Then little bird mom climbed onto the tree Then, the little bird mom climbed onto the tree.
	Predicate adjective	小猫 真 生气。 The cat so angry The cat is so angry.
Aspect Marker	Progressive	小狗 在 捉 老鼠。 The dog ProM catch mouse The dog is catching the mouse.
	Perfective	然后 小 男孩 拣回 了 自己的 球。 Then little boy got back PerM his ball Then the little boy got his ball back.
Classifier	-	小猫 吃 了 一 条 鱼。 The cat ate PerM one CL fish The cat at a fish.
Passive	-	它的 腿 被 狐狸 吃了。 Its leg PASS fox eat SFP Its leg was eaten by the fox.
Relative Clause	-	看到 掉 到 河 里 的 球 了。 Saw dropped RP river LP ball SFP Saw the ball that dropped in the water.

Note: ProM=progressive marker; PerM=perfective marker; CL=classifier; PASS=passive marker; SFP=sentence final particle; RP=resultative particle; LP=linking particle.